

(19)



Europäisches Patentamt

European Patent Office

Office européen des brevets



(11)

EP 0 881 591 B1

(12)

EUROPEAN PATENT SPECIFICATION

(45) Date of publication and mention
of the grant of the patent:
10.09.2003 Bulletin 2003/37

(51) Int Cl.7: **G06K 9/20**

(21) Application number: **98303390.3**

(22) Date of filing: **30.04.1998**

(54) Ordering groups of text in an image

Ordnen von Textgruppen in einem Bild

Mise en ordre de groupes de texte dans une image

(84) Designated Contracting States:
DE FR GB

(30) Priority: **29.05.1997 US 864993**
30.12.1997 US 898

(43) Date of publication of application:
02.12.1998 Bulletin 1998/49

(73) Proprietor: **Adobe Systems, Inc.**
San Jose, California 95110 (US)

(72) Inventor: **Stolin, Jacob**
Cupertino, California 95014 (US)

(74) Representative: **Wombwell, Francis et al**
Potts, Kerr & Co.
15, Hamilton Square
Birkenhead Merseyside CH41 6BR (GB)

(56) References cited:

- **ITO S ET AL: "Field segmentation and classification in document image" PROCEEDINGS OF THE 6TH INTERNATIONAL CONFERENCE ON PATTERN RECOGNITION, MUNICH, WEST GERMANY, 19-22 OCT. 1982, pages 492-495 vol.1, XP002076624 1982, New York, NY, USA, IEEE, USA**
- **BALESTRI M ET AL: "A METHOD FOR THE CORRECT ORDERING OF TYPEWRITTEN LINES" SIGNAL PROCESSING: THEORIES AND APPLICATIONS, GRENOBLE, SEPT. 5 - 8, 1988, vol. VOL. 3, no. CONF. 4, 5 September 1988, pages 1609-1611, XP000118027 LACOUME J L; CHEHIKIAN A; MARTIN N; MALBOS J**

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

Description

[0001] Paper documents can be scanned and stored as images in a computer. Text recognition techniques, such as optical character recognition (OCR), can then be used to convert text in these images to a computer-editable format, such as ASCII characters. Scanned images can contain text organized in multiple, distinct blocks (e.g., multiple columns of text, headlines, captions, footnotes, footers). The text blocks may further be separated by relatively large areas of blank space and graphical objects (lines, pictures, and so forth). Text can also be surrounded by a frame or contain insets, which further separate the text blocks into blocks. Although a person reading the page may be able to recognize the proper order of the text blocks in the image, it may be difficult for an OCR program to identify the text (by discarding the non-text components such as blank spaces and graphical objects) and then group the text into the proper reading order.

[0002] Examples of prior art arrangements are discussed in ITO et al.: 'Field segmentation and classification in document image' Proceedings of the 6th Int. Conf. On Pattern Recognition, Munich, Germany, 19-22 Oct. 1982, pages 492-495 vol.1, 1982, IEEE New York, NY, USA and BALESTRI et al.: 'A method for the correct ordering of typewritten lines' Signal Processing: Theories and Applications, Grenoble, Sept. 5 - 8, 1988, vol.3, no. Conf. 4,5 September 1988, pages 1609-1611.

SUMMARY

[0003] According to the present invention there is provided a computer-implemented method for ordering text in an image stored in a computer, the text being grouped in multiple blocks, the method comprising:

- grouping the text in multiple regions;
- representing the text regions as a graph having vertices and edges;
- defining each text region as vertex in the graph;
- defining edges between the vertices in the graph;
- assigning weights to the edges; and
- calculating a shortest Hamiltonian path through the vertices according to the edge weights; and
- ordering the text regions according to the order defined by the calculated shortest Hamiltonian path.

[0004] Also according to the present invention there is provided a program residing on a computer-readable medium for ordering text in an image stored in a computer, the program comprising instructions for causing the computer to:

- group the text in multiple regions;
- represent the text regions as a graph having vertices and edges;
- define each text region as a vertex in the graph;
- define edges between the vertices in the graph;
- assign weights to the edges in the graph; and
- calculate a shortest Hamiltonian path through the vertices according to the edge weights; and
- order the text regions according to the order defined by the calculated shortest Hamiltonian path.

[0005] Further according to the present invention there is provided apparatus for recognizing text in an image, comprising: a storage medium to store the image; and

- a processor operatively coupled to the storage medium and configured to:
- group the text in multiple regions;
- represent the text regions in a graph having vertices and edges;
- define each text region as a vertex in the graph;
- define edges between the vertices in the graph;
- assign weights to the edges in the graph; and
- calculate a shortest Hamiltonian path through the vertices according to the edge weights; and
- order the text regions according to the order defined by the calculated shortest Hamiltonian path.

[0006] Likewise according to the present invention there is provided a method implemented in a computer for ordering text in an image stored in the computer, the method comprising:

- identifying a set of text blocks;
- separating the set of text blocks into independent subsets of text blocks;

representing the text blocks as vertices in a graph in each subset;
 defining directed edges between vertices in each subset;
 assigning weights to the directed edges;
 calculating a shortest Hamiltonian path through the graph in each subset according to the edge weights;
 ordering the text blocks in each subset according to the order defined by the calculated shortest Hamiltonian path;
 and
 concatenating the ordering of text blocks in the subsets into a final order.

[0007] The invention has one or more of the following advantages. The proper order of multiple, distinct blocks of text in a captured image can be determined reliably by a text capture program.

[0008] Other features and advantages of the invention will become apparent from the following description and from the claims.

BRIEF DESCRIPTION OF THE DRAWINGS

[0009]

Fig. 1 is a flow diagram of a text capturing and ordering program in accordance with the present invention.

Fig. 2a is a diagram of text in an image separated into blocks.

Fig. 2b is a diagram of vertices representing the text blocks of Fig. 2a.

Fig. 3 is a graph containing the vertices corresponding to the text blocks of Fig. 2a along with oriented pairs of edges adjacent any two vertices.

Fig. 4 is a diagram of an optimal Hamiltonian path through the vertices in the graph of Fig. 3.

Fig. 5 is a diagram of a text block.

Fig. 6 is a flow diagram of a process for separating a page into independent parts.

Figs. 7, 8, and 9 are diagrams of text blocks in page parts.

Fig. 10 is a block diagram of a computer system.

DETAILED DESCRIPTION

[0010] Referring to Fig. 1, a computer implemented text capturing program is described that can reliably identify the proper order of text grouped into multiple, distinct blocks in an image. The program first captures and stores an image (step 102). Next, the program identifies the text blocks in the page based on conventional page layout analyses (step 104). For example, the image can be represented as density histograms, with very dense regions indicating non-text objects, such as graphical objects, and very sparse regions indicating gaps. Alternatively, the identification of text blocks can also be based on such factors as the proximity of the text blocks to each other, font size, and the existence of space separators and blocks of graphical objects. Thus, for example, although text characters in a page may be horizontally aligned, they may be separated by a wide gap, indicating that the characters are located in two different columns. In addition, the section heading for the page of text may have a different, larger font than the remaining text. The text characters may also be separated by graphical objects interspersed throughout the page.

[0011] After the text blocks have been identified in the image, the program separates the text blocks into independent subsets or parts of the page, if possible (step 106). Many pages can be divided into smaller parts that are divided by certain types of separators. These independent parts can be processed separately by the program, thereby reducing the complexity of finding the order of the text blocks in a page. Steps 108-116 in Fig. 1 are performed separately for each identified independent part of the page.

[0012] To further reduce the complexity of finding the order of the text blocks in each part of the page, the program next combines text blocks where possible (step 108). Often there is only one way to order two or more text blocks. In such cases, the blocks can be combined into a new single text region.

[0013] In the exemplary part 200 of a page in Fig. 2a, the darkened boxes 202 and 204 correspond to non-text objects, such as graphical objects. Further, a vertical divider line 206 separates text. In this image, the identified text blocks are labeled as text blocks 1-8. In each page part, the program then designates each text region (a region can be one text block or a group of combined text blocks) as a vertex of a graph (step 110). In Fig. 2a, text blocks 1 and 2 can be combined into one text region 12 and text blocks 6 and 7 can be combined into one text region 67. Thus, in Fig. 2b, vertices V_{12} , V_3 , V_4 , V_5 , V_{67} , and V_8 are designated for the text regions in Fig. 2a. The positions of the vertices are not necessarily geometrically related to the locations of the text blocks 1-8 in the image 200.

[0014] Next, the program defines directed edges (V_i, V_j) and (V_j, V_i) for each pair of vertices V_i and V_j (step 112). A pair of directed or oriented edges is defined between any two vertices because of the possibility that, as between any two text regions, one text region may come before the other text region. The vertices V_{12} , V_3 , V_4 , V_5 , V_{67} , and V_8 along

with the directed edges between each of the vertices define a directed or oriented graph G , as shown in Fig. 3.

[0015] The relationships between the vertices V are then defined (step 114) by assigning edge lengths (or weights) to the directed edges (V_i, V_j) and (V_j, V_i) , based on a number of factors. These factors include the distance between any two text blocks, the characteristics (e.g., number of lines, font size, spacing) of two text blocks, and the existence of separators (such as empty space or other non-text objects) between the text block pairs. The edge lengths are based on the likelihood that: one vertex V_i comes before its adjacent vertex V_j . The higher the likelihood that text region i comes before text region j , the smaller the weight of edge (V_i, V_j) , and vice versa.

[0016] Thus, for example, in Fig. 3, the weight assigned to the edge (V_{12}, V_3) is much smaller than the weight assigned the edge (V_3, V_{12}) because it is much more likely that text region 12 comes before text region 3.

[0017] Next, using the weights determined for the edges of the graph, the program finds an optimal Hamiltonian path through the vertices $V_{12}, V_3, V_4, V_5, V_{67}, V_8$ by using brute force (for small graphs) or conventional heuristic or approximate methods that solve a traveling salesman problem (step 116). An identified optimal Hamiltonian path is shown in Fig. 4, with the path starting at V_{12} and continuing to vertices V_3, V_4, V_5, V_{67} , and V_8 successively. Next, the program combines the partial orders found for the corresponding parts of the page into a final order π (step 118).

[0018] The following mathematical model is defined to perform the text ordering process. Referring to Fig. 5, for a text region A with coordinates (T, B, L, R) in a two-dimensional X - Y space, let

$$\begin{aligned} Top(A) &= T, \\ Bot(A) &= B, \\ Lft(A) &= L, \\ Rgt(A) &= R, \end{aligned} \quad (\text{Eq. 1})$$

and

$$\begin{aligned} CntrX(A) &= (L + R) / 2, \\ CntrY(A) &= (T + B) / 2, \end{aligned} \quad (\text{Eq. 2})$$

where L and R are on the X axis and T and B are on the Y axis. The distance between any two text regions $A1$ and $A2$ is defined as

$$|A1, A2| = |CntrX(A1) - CntrX(A2)| + |CntrY(A1) - CntrY(A2)| \quad (\text{Eq. 3})$$

[0019] For each pair of text regions A_i, A_j , a precedence function $f(A_i, A_j)$ is constructed so that, the more likely A_i precedes A_j , the smaller is the value of $f(A_i, A_j)$. For K text regions, the precedence functions $f(A_i, A_j)$, $i=1-K, j=1-K, i \neq j$, are calculated, which are used to calculate the edge lengths or weights between vertices.

[0020] However, before the precedence functions $f(A_i, A_j)$ are constructed, the complexity of the problem is reduced (1) by separating the page into different parts; and (2) by combining text blocks into regions, where possible.

[0021] Referring to Fig. 6, the step 106 of splitting the page into multiple parts is described. A page can be split into independent parts by applying the following recursive algorithm. The program creates a set SP of page parts P_i and initiates it by defining the whole page as the element of the set SP (step 300). For each element in SP , the program looks for a splitting separator going through the existing element (step 302), where a splitting separator may be defined as any non-text region except a thin vertical line, which might be a column separator. If no splitting separator is found (step 304), the process is stopped. If a splitting separator is found, the current element is divided into 2 new sub-elements by splitting it along the selected separator (step 306). The current element is replaced by two new sub-elements and steps 302-306 are repeated. The sub-elements form the parts P_i of the page. Thus $SP = \{P_1, P_2, \dots, P_n\}$, and the text ordering process is performed independently on each part P_i , with the results for each part combined at the end to determine the final order π .

[0022] To reduce the complexity in each page part, two or more text blocks or regions can be combined (step 108) if they are "horizontally connected" or "vertically connected." Two text regions $A1, A2$ are called horizontally connected (see Fig. 7) if the following conditions are true:

(1) A1 and A2 are horizontally aligned, that is, $\max(\text{Top}(A1), \text{Top}(A2)) < \min(\text{CntrY}(A1), \text{CntrY}(A2))$,
and
 $\min(\text{Bot}(A1), \text{Bot}(A2)) > \max(\text{CntrY}(A1), \text{CntrY}(A2))$;

(2) no other region overlaps a common bounding box of A1 and A2;

(3) A1 and A2 are blocked at the top, which means there are no regions above A1 and A2 or the nearest region A3 above A1 and A2 is a blocking region, that is,

$\text{Lft}(A3) \leq \min(\text{Lft}(A1), \text{Lft}(A2))$, and
 $\text{Rgt}(A3) \geq \max(\text{Rgt}(A1), \text{Rgt}(A2))$; and

(4) A1 and A2 are blocked at the bottom, that is, there are no regions below A1 and A2 or the nearest region A3 below A1 and A2 is a blocking region.

If the regions A1, A2 are horizontally connected, their partial order is from the left to the right (from A1 to A2 in Fig. 7).

[0023] Two text regions A1, A2 are vertically connected (see Fig. 8), if the following conditions are true:

(1) A1 and A2 are vertically aligned; that is, $\max(\text{Lft}(A1), \text{Lft}(A2)) < \min(\text{CntrX}(A1), \text{CntrX}(A2))$, and
 $\min(\text{Rgt}(A1), \text{Rgt}(A2)) > \max(\text{CntrX}(A1), \text{CntrX}(A2))$;

(2) no other region overlaps their common bounding box;

(3) A1 and A2 are blocked at the left, that is, there are no regions at the left or the nearest region A3 at the left is a blocking region, that is,

$\text{Top}(A3) \leq \min(\text{Top}(A1), \text{Top}(A2))$, and
 $\text{Bot}(A3) \geq \max(\text{Bot}(A1), \text{Bot}(A2))$; and

(4) A1 and A2 are blocked at the right; that is, there are no regions at the right or the nearest region A3 at the right is a blocking region.

[0024] If the regions A1, A2 are vertically connected, their partial order is from the top to the bottom (from A1 to A2 in Fig. 8).

[0025] If a pair of (horizontally or vertically) connected regions A1, A2 is found, the regions are combined into single text region A12. The bounding box of the new region is the smallest rectangle covering both A1 and A2 so that

$$\begin{aligned} \text{Top}(A) &= \min(\text{Top}(A1), \text{Top}(A2)), \\ \text{Lft}(A) &= \min(\text{Lft}(A1), \text{Lft}(A2)), \\ \text{Rgt}(A) &= \max(\text{Rgt}(A1), \text{Rgt}(A2)), \text{ and} \\ \text{Bot}(A) &= \max(\text{Bot}(A1), \text{Bot}(A2)). \end{aligned} \quad (\text{Eq. 4})$$

Other parameters (such as font size and spacing) for the combined region could be transferred from the bigger of regions A1, A2.

[0026] The combining process can be repeated until no more connected regions are found.

[0027] In some cases, the order of the text regions in a page part can be identified just by consecutively combining connected text regions. For example, in the page layout shown in Fig. 9, the solution could be found by combining the text regions as follows:

combine A5 and A6 into A56;
combine A56 and A7 into A567;
combine A2 and A567 into A2567;
combine A2567 and A8 into A25678;
combine A1 and A25678 into A125678;
combine A125678 and A3 into A1256783;
combine A1256783 and A4 into A12567834;

The resultant order of the uncombined text blocks is then A1, A2, A5, A6, A7, A8, A3, and A4.

[0028] After all connected text regions are combined, if more than one text region remains in a page part, the order of the regions is determined by solving for the optimal Hamiltonian path of a graph G containing vertices V representing

the uncombined text regions.

[0029] For K text regions, this is accomplished by first constructing precedence functions $f(A_i, A_j)$, $i \neq j$, $i = 1-K$, $j = 1-K$, for all the text regions. The precedence functions are used to assign lengths or weights to edges F_{ij} between vertices V_i and V_j .

[0030] The precedence function is defined as

$$f(A_i, A_j) = K_{loc}(A_i, A_j) + K_{dif}(A_i, A_j) + K_{sep}(A_i, A_j), \quad \text{Eq. 5}$$

where K_{loc} evaluates the relative locations of the two text regions A_i and A_j ; K_{dif} evaluates the similarity (in number of lines, font size, and spacing) of text regions; and K_{sep} reflects the contribution to the function f due to the existence of a non-text separating region, if any, between A_i and A_j . How K_{loc} , K_{dif} , and K_{sep} are derived is described below.

[0031] A graph G associated with a page part is defined as follows: G is a directed graph with K vertices V_1, V_2, \dots, V_K ; each pair of vertices V_i, V_j , $i \neq j$ is connected by a directed edge E_{ij} ; a non-negative number $W(E_{ij})$ (referred to as the weight or length of the edge E_{ij}) is assigned to each edge E_{ij} :

$$W(E_{ij}) = f(A_i, A_j), \quad \text{(Eq. 6)}$$

where f is the precedence function defined by Eq. 5.

[0032] For a given order π (which is a permutation of numbers 1, 2, ..., k), a Hamiltonian path $P(\pi)$ in the graph G is an ordered set of vertices

$$P(\pi) = \{V_{\pi(1)}, V_{\pi(2)}, \dots, V_{\pi(k)}\}. \quad \text{(Eq. 7)}$$

[0033] The length of the path $P(\pi)$ is defined as

$$L(\pi) = W(E_{\pi(1)\pi(2)}) + W(E_{\pi(2)\pi(3)}) + \dots + W(E_{\pi(k-1)\pi(k)}). \quad \text{(Eq. 8)}$$

[0034] The shortest Hamiltonian path is a Hamiltonian path with the minimal value $L(\pi)$.

[0035] Each Hamiltonian path $P(\pi)$ in the associated graph G defines an order of text regions in the corresponding page part. As it follows from the definition of the precedence function f , the shorter the Hamiltonian path $P(\pi)$, the greater the likelihood that π is the proper logical order of text regions. Therefore, the shortest Hamiltonian path $P(\pi)$ in the graph G provides the solution for finding the order π of text blocks in a page part.

[0036] To find the shortest Hamiltonian path, the standard method of reducing it to the traveling salesman problem can be used. First, an additional vertex V_0 is added to the graph G, with the vertex V_0 connected to each vertex V_j by edges E_{0j} and E_{j0} . The length of each edge E_{0j} and E_{j0} is 0, i.e., $W(E_{0j}) = W(E_{j0}) = 0$. Next, a shortest ordered cycle C in the graph G is calculated by applying a standard algorithm for solving traveling salesman problems. The shortest Hamiltonian path is then extracted from the cycle C by removing the additional vertex V_0 from the cycle.

[0037] Once the logical orders π_j for independent parts P_1, \dots, P_k in the page have been identified, the paths π_j , $j=1-n$, are concatenated:

$$\pi = \{\pi_1, \pi_2, \dots, \pi_n\} \quad \text{Eq. 9}$$

[0038] Since the parts P are independent, it does not matter how the orders π are concatenated. However, an alternative is to sort the parts P in increasing order of y and then x where (x,y) is the top, left corner of each page part.

[0039] In the concatenated order π , combined text blocks in text regions are separated out and placed in proper order in a final order π' . Thus, for example, if π is {12, 3, 7, 56, 4}, it is modified to $\pi' = \{1, 2, 3, 7, 5, 6, 4\}$ to define the order of text blocks A1-A7.

[0040] The final order π' thus provides the solution for the problem of identifying the order of text blocks in a page.

[0041] As set forth in Eq. 5, the precedence functions $f(A_i, A_j)$, $i \neq j$, $i = 1-k$, $j = 1-k$ are calculated based on values $K_{loc}(A_i, A_j)$, $K_{diff}(A_i, A_j)$, and $K_{sep}(A_i, A_j)$ for k text regions.

[0042] A row preference or column preference can be selected. If row preference is selected, then text region ordering favors ordering in the X direction. If column preference is selected, the text region ordering favors ordering in the Y direction. For regions $A1=(T1, B1, L1, R1)$ and $A2=(T2, B2, L2, R2)$, the component K_{loc} , which has a value that is dependent upon the relative coordinates of regions A1 and A2, is calculated differently for row and column preferences. Since K_{loc} is dependent on the relative locations of A1 and A2, $K_{loc}(A1, A2)$ is calculated differently than $K_{loc}(A2, A1)$, with $K_{loc}(A1, A2)$ used to calculate $f(A1, A2)$ and $K_{loc}(A2, A1)$ used to calculate $f(A2, A1)$. Generally, because one text region will come before the other region, $f(A1, A2)$ is usually not equal to $f(A2, A1)$ due to the differences in calculating $K_{loc}(A1, A2)$ and $K_{loc}(A2, A1)$. The calculation of $K_{loc}(A1, A2)$ or $K_{loc}(A2, A1)$ is set forth below for three possible cases. Since by definition two separate regions do not overlap each other, the case where $\min(R1, R2) \geq \max(L1, L2)$ and $\min(B1, B2) \geq \max(T1, T2)$ is not possible and thus not considered.

[0043] In a first case, the text regions A1 and A2 do not overlap either in the X or Y axis and A1 is to the left of and below A2; that is, $R1 < L2$ and $T1 > B2$.

[0044] In this case, if column preference is selected, the value K_{loc} is defined as

$$K_{loc}(A1, A2) = Q1 * |CntrX(A1) - CntrX(A2)|, \quad (\text{Eq. 10})$$

where Q1 is a tunable parameter with a default value of 1; and

$$K_{loc}(A2, A1) = Q2 * |A1, A2|, \quad (\text{Eq. 11})$$

where Q2 is a tunable parameter with default value 2. Thus $K_{loc}(A1, A2)$ has a smaller value than $K_{loc}(A2, A1)$, which tends to favor A1 over A2, which is consistent with the selected column preference.

[0045] In the first case, if row preference is selected, the value K_{loc} is defined as:

$$K_{loc}(A1, A2) = Q3 * |A1, A2|, \quad (\text{Eq. 12})$$

where Q3 is a tunable parameter with default value 4, and

$$K_{loc}(A2, A1) = Q4 * |A1, A2|, \quad (\text{Eq. 13})$$

where Q4 is a tunable with default value 1. $K_{loc}(A1, A2)$ has a larger value than $K_{loc}(A2, A1)$, which tends to favor A2 over A1 if row preference is selected in the first case.

[0046] In a second case, the boundaries of the text regions A1 and A2 do not overlap in the X axis but overlap in the Y axis and the region A1 is to the left of region A2; that is, $R1 < L2$ and $T1 \leq B2$.

[0047] In this case, if column preference is selected, the value K_{loc} is defined as:

$$K_{loc}(A1, A2) = Q1 * |CntrX(A1) - CntrX(A2)|, \quad (\text{Eq. 14})$$

and

$$K_{loc}(A1, A2) = Q1 * |CntrX(A1) - CntrX(A2)|, \quad (\text{Eq. 14})$$

and

$$K_{loc}(A2,A1) = M1, \quad (\text{Eq. 15})$$

where M1 is a large value, which can be set to

$$M1 = 10 * \max_{ij} |Ai, Aj|, \quad (\text{Eq. 16})$$

[0048] M1 is thus defined as ten times the maximum possible distance between any two text regions in a considered part of the page. Generally, this heavily favors A1 over A2 in the calculation of K_{loc} .

[0049] In the second case, if row preference is selected, the value K_{loc} is defined as:

$$K_{loc}(A1,A2) = Q4 * |A1,A2|, \quad (\text{Eq. 17})$$

and

$$K_{loc}(A2,A1) = M1, \quad (\text{Eq. 18})$$

where Q4 is tunable with a default value of 1. Again, A1 is generally heavily favored over A2.

[0050] In a third case, the boundaries of the text regions do not overlap in the Y axis but overlap in the X axis and A1 is located above A2; that is, $B1 < T2$ and $\min(R1,R2) \geq \max(L1,L2)$.

[0051] In this case, for both column and row preferences, the value K_{loc} is defined as:

$$K_{loc}(A1,A2) = Q5 * |A1,A2|, \quad (\text{Eq. 19})$$

and

$$K_{loc}(A2,A1) = M1, \quad (\text{Eq. 20})$$

where Q5 is a tunable parameter with a default value 1. Those calculations generally heavily favor A1 over A2.

[0052] The function $K_{dif}(A1,A2)$ is defined as

$$K_{dif}(A1,A2) = Q6 * (m_1 + m_2) * (|s_1 - s_2| + |l_1 - l_2|), \quad (\text{Eq. 21})$$

where m_1 is the number of text lines in region Ai , s_i is the text point size for region Ai , l_i is the distance between consecutive lines in Ai , and Q6 is a tunable parameter (default value is 10). In effect, s_i and l_i represent the height (in the Y direction) of a line. $K_{dif}(A2,A1)$ is equal to $K_{dif}(A1,A2)$.

[0053] The function $K_{sep}(A1,A2)$ is defined relative to a separator (non-text region) B and is calculated in terms of a horizontal extrusion parameter $E_{hor}(A,B)$ and a vertical extrusion parameter $E_{vert}(A,B)$. For a text region A and separator B, the horizontal extrusion parameter $E_{hor}(A,B)$ is defined as

$$E_{hor}(A,B) = \max[Lft(A) - Lft(B), Rgt(B) - Rgt(A)] / [Rgt(A) - Lft(A)],$$

if $[Lft(B) < Lft(A) \text{ and } Rgt(B) > Lft(A)]$

or $[Rgt(B) > Rgt(A) \text{ and } Lft(B) < Rgt(A)]$;

0, otherwise; (Eq. 22)

[0054] Thus, $E_{hor}(A,B)$ is greater than zero if either of the left or right edges of the text region A falls within the range defined by the left and right edges of the separator B.

[0055] Similarly, a vertical extrusion parameter $E_{vert}(A,B)$ is defined as

$$E_{vert}(A,B) = \max [Top(A) - Top(B), Bot(B) - Bot(A)] /$$

$$[Bot(A) - Top(A)],$$

if $[Top(B) < Top(A) \text{ and } Bot(B) > Top(A)]$

or $[Bot(B) > Bot(A) \text{ and } Top(B) < Bot(A)]$;

0, otherwise. (Eq. 23)

[0056] The function $K_{sep}(A1,A2)$ is defined as follows for the following two possible cases.

[0057] In a first case, the text regions A1, A2 are vertically disjoint; that is, the regions A1,A2 do not overlap in the Y axis, defined by $\min(Bot(A1), Bot(A2)) < \max(Top(A1), Top(A2))$. In this case,

$$K_{sep}(A1, A2) = Q7 * \sum_B (E_{hor}(A1, B) + E_{hor}(A2, B)) , \quad (\text{Eq. 24})$$

where Q7 is a tunable parameter (default can be 10) and the sum \sum_B includes all separators between A1 and A2; i.e.,

$Top(B) > \min(Bot(A1), Bot(A2))$, and
 $Bot(B) < \max(Top(A1), Top(A2))$.

[0058] In a second case, the text regions are horizontally disjoint; that is, the regions A1,A2 do not overlap on the X axis as defined by $\min(Rgt(A1), Rgt(A2)) < \max(Lft(A1), Lft(A2))$. In this case,

$$K_{sep}(A1, A2) = Q7 * \sum_B (E_{vert}(A1, B) + E_{vert}(A2, B)) , \quad (\text{Eq. 25})$$

where the sum \sum_B includes all separators between A1 and A2; i.e.,

$Lft(B) > \min(Rgt(A1), Rgt(A2))$, and
 $Rgt(B) < \max(Lft(A1), Lft(A2))$.

[0059] Once K_{loc} , K_{dif} , and K_{loc} are calculated for all combinations of A1, A2, ..., Ak, the precedence functions $f(Ai, Aj)$, $i \neq j$, $i=1-k$, $j=1-k$, can be constructed and used in finding the lengths of different permutations of paths $P(\pi)$ to identify the shortest Hamiltonian path $P(\pi)$.

[0060] Referring to Fig. 10, the text capturing and ordering program may be implemented in digital electronic circuitry or in computer hardware, firmware, software, or in combinations of them, such as in a computer system. The computer system includes a central processing unit (CPU) 502 connected to an internal system bus 504. The storage media in the computer system include a main memory 506 (which can be implemented with dynamic random access memory devices), a hard disk drive 508 for mass storage, and a non-volatile memory (NVRAM) 510. The main memory 506 and NVRAM 510 are connected to the bus 504, and the hard disk drive 508 is coupled to the bus 504 through a hard disk drive controller 512.

[0061] Apparatus of the invention may be implemented in a computer program product tangibly embodied in a machine-readable storage device (such as the hard disk drive 508, main memory 506, or NVRAM 510) for execution by the CPU 502. Suitable processors include, by way of example, both general and special purpose microprocessors.

Generally, a processor will receive instructions and data from the read-only memory 510 and/or the main memory 506. Storage devices suitable for tangibly embodying computer programming instructions include all forms of non-volatile memory, including by way of example semiconductor memory devices, such as EPROM, EEPROM, and flash memory devices; magnetic disks such as the internal hard disk drive 508 and removable disks and diskettes 528 connected
 5 through a controller 526; magneto-optical disks; and CD-ROM disks. Any of the foregoing may be supplemented by, or incorporated in, specially-designed ASICs (application-specific integrated circuits).

[0062] The computer system further includes an input-output (I/O) controller 514 connected to the bus 504 and which provides a keyboard interface 516 for connection to an external keyboard, a mouse interface 518 for connection to an external mouse or other pointer device, and a parallel port interface 520 for connection to a printer. In addition, the bus
 10 504 is connected to a video controller 522 which couples to an external computer monitor or display 524. Data associated with an image for display on a computer monitor 524 are provided over the system bus 504 by application programs to the video controller 522 through the operating system and the appropriate device driver.

[0063] Other embodiments are within the scope of the following claims. For example, the order of the steps of the invention may be changed by those skilled in the art and still achieve desirable results. Different techniques can be
 15 used to identify an optimal path between vertices of a graph representing text blocks or regions in an image. Although specific equations and parameters have been disclosed to determine variables used in finding an optimal order of text blocks or regions, such equations and parameters can be changed.

20 Claims

1. A computer-implemented method for ordering text in an image stored in a computer, the text being grouped in multiple blocks, the method comprising:

25 grouping the text in multiple regions;
 representing the text regions as a graph having vertices and edges (110);
 defining each text region as vertex in the graph;
 defining edges between the vertices in the graph (112);

30 said method being **characterized by** the steps of:

assigning weights to the edges (114); and
 calculating a shortest Hamiltonian path through the vertices according to the edge weights (116); and
 ordering the text regions according to the order defined by the calculated shortest Hamiltonian path (113).

2. The method of claim 1, wherein oriented pairs of edges are defined between any two vertices.

3. The method of claim 1, wherein the step of calculating a shortest Hamiltonian path (116) comprises:

40 adding a virtual vertex and virtual oriented edges to the graph;
 obtaining a shortest ordered cycle in the graph by solving a traveling salesman problem on the graph; and
 obtaining the shortest Hamiltonian path by removing the virtual vertex from the shortest ordered cycle.

4. The method of claim 1, wherein the weights assigned by the edges between the vertices are based on the distance
 45 between corresponding text regions.

5. The method of claim 1, wherein the weights assigned the edges between vertices are based on the text characteristic of the corresponding text regions.

- 50 6. The method of claim 5, wherein the text characteristics include font size and number of lines of text.

7. The method of claim 1, wherein the weights assigned the edges between the vertices are based on the existence of non-text separators between text region pairs.

- 55 8. The method of claim 8, wherein the separators include graphical objects.

9. The method of claim 1, further comprising:

identifying text blocks that can be combined; and
combining the text blocks into a text region.

10. The method of claim 9, wherein two text blocks can be combined if they are vertically connected.

11. The method of claim 9, wherein two text blocks can be combined if they are horizontally connected.

12. The method of claim 1, further comprising:

initially separating the image into independent parts, each part containing its own set of text regions; and
performing the actions of grouping, representing, defining, assigning, calculating, and ordering on the text
regions in each part independently.

13. The method of claim 12, wherein the image is separated by identifying predetermined types of non-text separators.

14. The method of claim 12, further comprising:

concatenating the ordering of the text regions identified for the different parts.

15. A program residing on a computer-readable medium for ordering text in an image stored in a computer, the program
comprising instructions for causing the computer to:

group the text in multiple regions;
represent the text regions as a graph having vertices and edges (110);
define each text region as a vertex in the graph;
define edges between the vertices in the graph (112);

said program being **characterized in that** it comprises instructions for causing said computer to:

assign weights to the edges in the graph (114); and
calculate a shortest Hamiltonian path through the vertices according to the edge weights (116); and
order the text regions according to the order defined by the calculated shortest Hamiltonian path (118).

16. The program of claim 15, wherein the weights assigned the edges between the vertices are based on the distance
between corresponding text blocks and the characteristics of each block.

17. The program of claim 15, wherein the weights assigned the edges between the vertices are based on the existence
of separators between the text block pairs.

18. The program of claim 15, wherein the program comprises instructions for further causing the computer to:

identify blocks of text that can be combined; and
combine the text blocks into a text region.

19. The program of claim 15, wherein the program comprises instructions for further causing the computer to:

initially separate the image into independent parts, wherein each part contains its own set of text regions; and
separately perform the actions of grouping, representing, defining, assigning, calculating, and ordering the
text regions in each independent part.

20. The program of claim 19, wherein the program comprises instructions for further causing the computer to concate-
nate the ordering of text regions identified for the independent parts.

21. Apparatus for recognizing text in an image, comprising: a storage medium to store the image; and

a processor operatively coupled to the storage medium and configured to:
group the text in multiple regions;
represent the text regions in a graph having vertices and edges (110);

define each text region as a vertex in the graph;
define edges between the vertices in the graph (112);

said processor being **characterized in that** it is configured to:

assign weights to the edges in the graph (114); and
calculate a shortest Hamiltonian path through the vertices according to the edge weights (116); and
order the text regions according to the order defined by the calculated shortest Hamiltonian path (118).

22. The apparatus of claim 21, wherein the weights assigned the edges between the vertices are based on the distance between any two text regions.

23. The apparatus of claim 21, wherein the weights assigned the edges between the vertices are further based on the existence of separators between the text block pairs.

24. A method implemented in a computer for ordering text in an image stored in the computer, the method comprising:

identifying a set of text blocks (104);
separating the set of text blocks into independent subsets of text blocks (108);
representing the text blocks as vertices in a graph in each subset (110);
defining directed edges between vertices in each subset (112);

said method being **characterized by** the steps of:

assigning weights to the directed edges (114);
calculating a shortest Hamiltonian path through the graph in each subset according to the edge weights (116);
ordering the text blocks in each subset according to the order defined by the calculated shortest Hamiltonian path (118); and
concatenating the ordering of text blocks in the subsets into a final order.

Patentansprüche

1. Ein Computer-implementiertes Verfahren zum Ordnen von Text in einem in einem Computer gespeicherten Bild, wobei der Text in mehrere Blöcke gruppiert ist, wobei das Verfahren umfaßt:

Gruppieren des Texts in mehrere Gebiete;
Darstellen der Textgebiete als Graph mit Knoten und Rändern (110) ;
Definieren jedes Textgebietes als Knoten in dem Graph;
Definieren von Rändern zwischen den Knoten in dem Graph (112);

wobei das Verfahren **gekennzeichnet ist durch** die Schritte:

Zuweisen von Gewichten zu den Rändern (114); und
Berechnen eines kürzesten Hamilton-Pfades **durch** die Knoten entsprechend den Rändergewichten (116);
und
Ordnen der Textgebiete gemäß der **durch** den berechneten kürzesten Hamilton-Pfad definierten Reihenfolge.

2. Das Verfahren nach Anspruch 1, wobei orientierte Paare von Rändern zwischen zwei beliebigen Knoten definiert werden.

3. Das Verfahren nach Anspruch 1, wobei der Schritt des Berechnens eines kürzesten Hamilton-Pfades (116) umfaßt:

Hinzufügen eines virtuellen Knotens und virtueller orientierter Ränder zu dem Graphen;
Gewinnen einer kürzesten geordneten Schleife in dem Graphen durch Lösen eines Handlungsreisenden-Problems an dem Graphen; und
Gewinnen des kürzesten Hamilton-Pfades durch Entfernen des virtuellen Knotens aus der kürzesten geordneten Schleife.

4. Das Verfahren nach Anspruch 1, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf dem Abstand zwischen zugehörigen Textgebieten basieren.
- 5 5. Das Verfahren nach Anspruch 1, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf den Textcharakteristiken der zugehörigen Textgebiete basieren.
6. Das Verfahren nach Anspruch 5, wobei die Textcharakteristiken eine Font-Größe und eine Anzahl von Zeilen des Textes enthalten.
- 10 7. Das Verfahren nach Anspruch 1, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf dem Vorhandensein von Nicht-Text-Trennelementen zwischen Textgebietpaaren basieren.
8. Das Verfahren nach Anspruch 8, wobei die Trennelemente grafische Objekte enthalten.
- 15 9. Das Verfahren nach Anspruch 1, ferner umfassend:

Identifizieren von Textblöcken, die kombiniert werden können; und
Kombinieren der Textblöcke zu einem Textgebiet.
- 20 10. Das Verfahren nach Anspruch 9, wobei zwei Textblöcke kombiniert werden können, wenn sie vertikal verbunden sind.
11. Das Verfahren nach Anspruch 9, wobei zwei Textblöcke kombiniert werden können, wenn sie horizontal verbunden sind.
- 25 12. Das Verfahren nach Anspruch 1, ferner umfassend:

anfängliches Trennen des Bildes in unabhängige Teile, wobei jeder Teil seine eigene Menge von Textgebieten enthält; und
30 unabhängiges Ausführen der Aktionen des Kopierens, Darstellen, Definierens, Zuweisens, Berechnens und Ordners an den Textgebieten in jedem Teil.
13. Das Verfahren nach Anspruch 12, wobei das Bild aufgeteilt wird, indem vorgegebene Arten von Nicht-Text-Trennelementen identifiziert werden.
- 35 14. Das Verfahren nach Anspruch 12, ferner umfassend:

Verketteten der Ordnung der Textgebiete, die für die verschiedenen Teile identifiziert sind.
- 40 15. Ein auf einem Computer-lesbaren-Medium befindliches Programm zum Ordnen von Text in einem in einem Computer gespeicherten Bild, wobei das Programm Befehle aufweist, die den Computer veranlassen,

den Text in mehrere Gebiete zu gruppieren;
45 die Textgebiete als Graphen mit Knoten und Rändern darzustellen (110);
jedes Textgebiet als Knoten in dem Graphen zu definieren;
Ränder zwischen den Knoten in dem Graphen zu definieren (112);
wobei das Programm **dadurch gekennzeichnet ist, daß** es Befehle aufweist, die den Computer veranlassen, Gewichte den Rändern in dem Graphen zuzuweisen (114); und
50 einen kürzesten Hamilton-Pfad durch die Knoten entsprechend den Rändergewichten zu berechnen (116); und
die Textgebiete gemäß der durch den berechneten kürzesten Hamilton-Pfad definierten Reihenfolge zu ordnen (118).
16. Das Programm nach Anspruch 15, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf dem Abstand zwischen zugehörigen Textblöcken und den Charakteristiken jedes Blocks basieren.
- 55 17. Das Programm nach Anspruch 15, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf dem Vorhandensein von Trennelementen zwischen den Textblockpaaren basieren.

18. Das Programm nach Anspruch 15, wobei das Programm Befehle aufweist, die den Computer ferner veranlassen:

Blöcke des Textes zu identifizieren, die kombiniert werden können; und
die Textblöcke in ein Textgebiet zu kombinieren.

19. Das Programm nach Anspruch 15, wobei das Programm Befehle aufweist, um den Computer ferner zu veranlassen:

anfänglich das Bild in unabhängige Teile aufzuteilen, wobei jeder Teil seine eigene Menge von Textgebieten enthält; und
die Aktionen des Gruppierens, Darstellens, Definierens, Zuweisens, Berechnens und Ordnen der Textgebiete in jedem Teil unabhängig auszuführen.

20. Das Programm nach Anspruch 19, wobei das Programm Befehle aufweist, um den Computer ferner zu veranlassen, die Ordnung der Textgebiete zu verketten, die für die unabhängigen Teile identifiziert sind.

21. Einrichtung zum Erkennen von Text in einem Bild, aufweisend:

ein Speichermedium zum Speichern des Bildes; und
einen betriebsmäßig mit dem Speichermedium gekoppelten Prozessor, der konfiguriert ist, um:

den Text in mehrere Gebiete zu gruppieren;
die Textgebiete in einem Graphen mit Knoten und Rändern darzustellen (110);
jedes Textgebiet als Knoten in dem Graphen zu definieren;
Ränder zwischen den Knoten in dem Graphen zu definieren (112);

wobei der Prozessor **dadurch gekennzeichnet ist, daß** er konfiguriert ist, um:

den Rändern in dem Graphen Gewichte zuzuweisen (114); und
einen kürzesten Hamilton-Pfad durch die Knoten in Übereinstimmung mit den Rändergewichten zu berechnen (116); und
die Textgebiete in Übereinstimmung mit der durch den berechneten kürzesten Hamilton-Pfad definierten Reihenfolge zu ordnen (118).

22. Die Einrichtung nach Anspruch 21, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte auf dem Abstand zwischen zwei beliebigen Textgebieten basieren.

23. Die Einrichtung nach Anspruch 21, wobei die den Rändern zwischen den Knoten zugewiesenen Gewichte ferner auf dem Vorhandensein von Trennelementen zwischen Textblockpaaren basieren.

24. Ein in einem Computer implementiertes Verfahren zum Ordnen von Text in einem in dem Computer gespeicherten Bild, wobei das Verfahren umfaßt:

Identifizieren einer Menge von Textblöcken (104);
Aufteilen der Menge von Textblöcken in unabhängige Untermengen von Textblöcken (106);
Darstellen der Textblöcke als Knoten in einem Graphen in jeder Untermenge (110);
Definieren gerichteter Ränder zwischen den Knoten in jeder Untermenge (112);

wobei das Verfahren **gekennzeichnet ist durch** die Schritte:

Zuweisen von Gewichten zu den gerichteten Rändern (114);
Berechnen eines kürzesten Hamilton-Pfads **durch** den Graphen in jeder Untermenge in Übereinstimmung mit den Rändergewichten (116);
Ordnen der Textblöcke in jeder Untermenge in Übereinstimmung mit der **durch** den berechneten kürzesten Hamilton-Pfad definierten Reihenfolge (118); und
Verketteten der Ordnung der Textblöcke in den Untermengen zu einer Endreihenfolge.

Revendications

1. Procédé implémenté sur ordinateur pour ordonner du texte dans une image mémorisée dans un ordinateur, le texte étant groupé en de multiples blocs, le procédé comprenant les étapes consistant à :

grouper le texte en de multiples régions ;
représenter les régions de texte comme un graphe ayant des sommets et des arêtes (110) ;
définir chaque région de texte comme sommet dans le graphe ;
définir des arêtes entre les sommets dans le graphe (112) ;
ledit procédé étant **caractérisé par** les étapes consistant à :

affecter des poids aux arêtes (114) ; et
calculer le chemin Hamiltonien le plus court à travers les sommets en conformité avec les poids des arêtes (116) ; et
ordonner les régions de texte en conformité avec l'ordre défini par le chemin Hamiltonien le plus court calculé (113).

2. Procédé selon la revendication 1, dans lequel des paires orientées d'arêtes sont définies entre deux sommets.

3. Procédé selon la revendication 1, dans lequel l'étape de calcul du chemin Hamiltonien le plus court (116) comprend les étapes consistant à :

ajouter un sommet virtuel et des arêtes orientées virtuelles au graphe ;
obtenir un cycle ordonné le plus court dans le graphe en résolvant un problème du représentant de commerce sur le graphe ; et
obtenir le chemin Hamiltonien le plus court en enlevant le sommet virtuel du cycle ordonné le plus court.

4. Procédé selon la revendication 1, dans lequel les poids affectés par les arêtes entre les sommets sont basés sur la distance entre les régions de texte correspondantes.

5. Procédé selon la revendication 1, dans lequel les poids affectés aux arêtes entre les sommets sont basés sur la caractéristique de texte des régions de texte correspondantes.

6. Procédé selon la revendication 5, dans lequel les caractéristiques de texte comprennent la taille de la police et le nombre de lignes de texte.

7. Procédé selon la revendication 1, dans lequel les poids affectés aux arêtes entre les sommets sont basés sur l'existence de séparateurs d'absence de texte entre les paires de régions de texte.

8. Procédé selon la revendication 8, dans lequel les séparateurs comprennent des objets graphiques.

9. Procédé selon la revendication 1, comprenant en outre les étapes consistant à :

identifier les blocs de texte qui peuvent être combinés ; et
combinaison des blocs de texte en une région de texte.

10. Procédé selon la revendication 9, dans lequel deux blocs de texte peuvent être combinés s'ils sont verticalement reliés.

11. Procédé selon la revendication 9, dans lequel deux blocs de texte peuvent être combinés s'ils sont reliés horizontalement.

12. Procédé selon la revendication 1, comprenant en outre les étapes consistant à :

séparer initialement l'image en parties indépendantes, chaque partie contenant son propre ensemble de régions de texte ; et
effectuer les actions de groupage, de représentation, de définition, d'affectation, de calcul et d'ordonnement sur les régions de texte dans chaque partie indépendamment.

13. Procédé selon la revendication 12, dans lequel l'image est séparée en identifiant des types prédéterminés de séparateur d'absence de texte.

14. Procédé selon la revendication 12, comprenant en outre les étapes consistant à :

concaténer l'ordonnancement des régions de texte identifiées pour les parties différentes.

15. Programme résidant sur un support lisible par ordinateur pour ordonner le texte dans une image mémorisée dans un ordinateur, le programme comprenant les instructions pour amener l'ordinateur à :

grouper le texte en de régions multiples ;

représenter les régions de texte comme un graphe ayant des sommets et des arêtes (110);

définir chaque région de texte comme sommet dans le graphe ;

définir des arêtes entre les sommets dans le graphe (112) ;

ledit programme étant **caractérisé en ce qu'il** comprend les instructions pour amener l'ordinateur à :

affecter des poids aux arêtes dans le graphe (114) ; et

calculer le chemin Hamiltonien le plus court à travers les sommets en conformité avec les poids des arêtes (116) ; et

ordonner les régions de texte en conformité avec l'ordre défini par le chemin Hamiltonien le plus court calculé (118).

16. Programme selon la revendication 15, dans lequel les poids affectés aux arêtes entre les sommets sont basés sur la distance entre les blocs de texte correspondants et les caractéristiques de chaque bloc.

17. Programme selon la revendication 15, dans lequel les poids affectés aux arêtes entre les sommets sont basés sur l'existence de séparateurs entre les paires de blocs de texte.

18. Programme selon la revendication 15, dans lequel le programme comprend des instructions pour amener en outre l'ordinateur à :

identifier les blocs de texte qui peuvent être combinés et
combiner les blocs de texte en une région de texte.

19. Programme selon la revendication 15, dans lequel le programme comprend des instructions pour amener en outre l'ordinateur à :

séparer initialement l'image en parties indépendantes, dans lequel chaque partie contient son propre ensemble de régions de texte ; et

effectuer séparément les actions de groupage, de représentation, de définition, d'affectation, de calcul, et d'ordonnancement des régions de texte dans chaque partie indépendante.

20. Programme selon la revendication 19, dans lequel le programme comprend des instructions pour amener en outre l'ordinateur à concaténer l'ordonnancement des régions de texte identifiées pour les parties indépendantes.

21. Appareil pour reconnaître du texte dans une image, comprenant : un support de stockage pour stocker l'image ; et
un processeur couplé fonctionnellement au support de stockage et configuré pour :

grouper le texte en des régions multiples ;

représenter les régions de texte en un graphe ayant des sommets et des arêtes (110) ;

définir chaque région de texte comme un sommet dans le graphe ;

définir des arêtes entre les sommets dans le graphe (112) ;

ledit processeur étant **caractérisé en ce qu'il** est configuré pour :

affecter des poids aux arêtes dans le graphe (114) ; et

calculer le chemin Hamiltonien le plus court à travers les sommets en conformité avec les poids des arêtes (116) ; et

ordonner les régions de texte en conformité avec l'ordre défini par le chemin Hamiltonien le plus court calculé (118).

5 22. Appareil selon la revendication 21, dans lequel les poids affectés aux arêtes entre les sommets sont basés sur la distance entre deux régions quelconques de texte.

23. Appareil selon la revendication 21, dans lequel les poids affectés aux arêtes entre les sommets sont de plus basés sur l'existence de séparateurs entre les paires de blocs de texte.

10 24. Procédé implémenté sur ordinateur pour ordonner le texte dans une image mémorisée dans l'ordinateur, le procédé comprenant les étapes consistant à :

15 identifier un ensemble de blocs de texte (104) ;

 séparer l'ensemble de blocs de texte en sous-ensembles indépendants de blocs de texte (106);

 représenter les blocs de texte comme sommets dans un graphe dans chaque sous-ensemble (110);

 définir des arêtes dirigées entre les sommets dans chaque sous-ensemble (112);

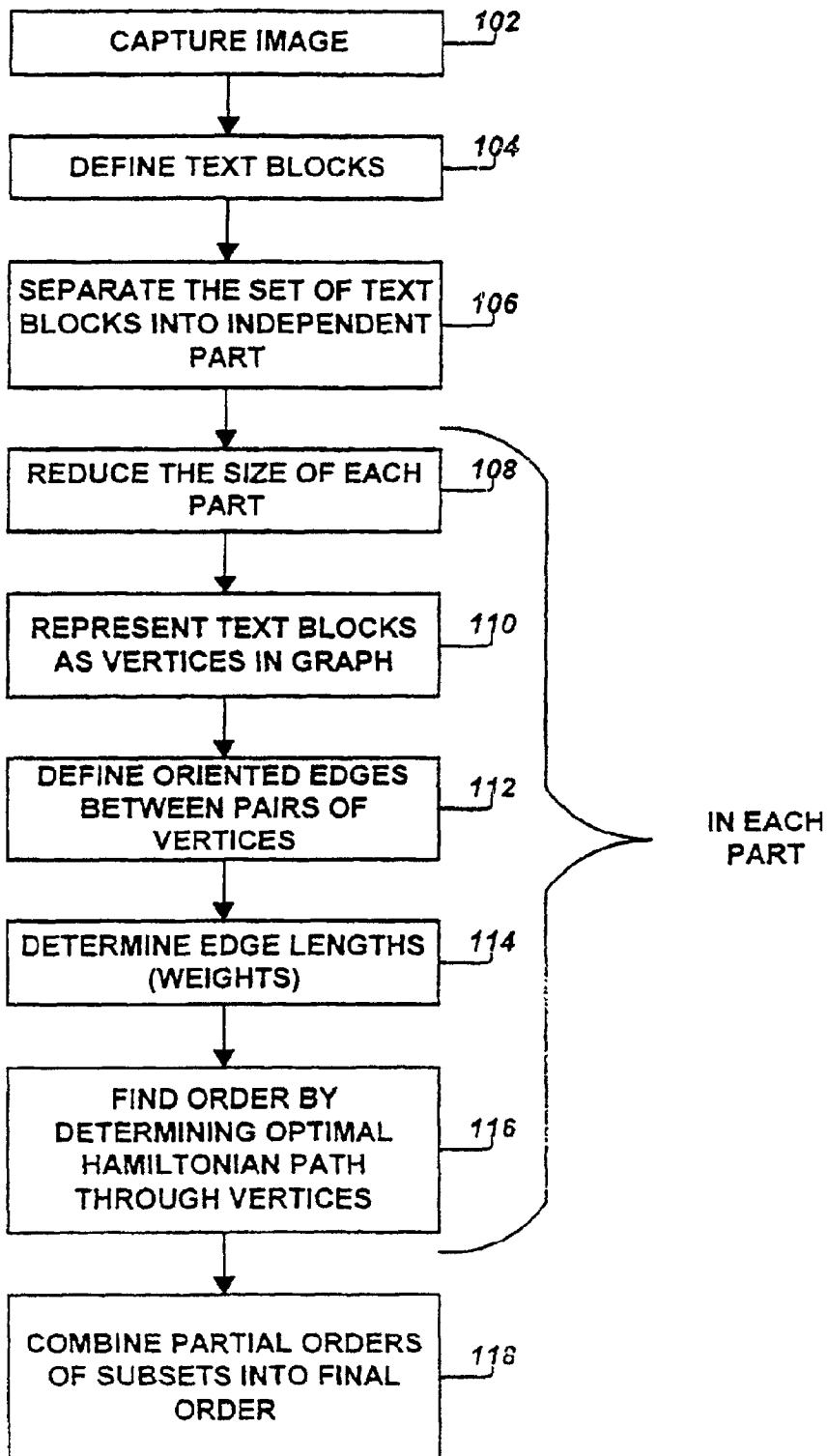
 ledit procédé étant **caractérisé par** les étapes consistant à :

20 affecter des poids aux arêtes dirigées (114) ;

 calculer le chemin Hamiltonien le plus court à travers le graphe dans chaque sous-ensemble en conformité avec les poids des arêtes (116) ;

 ordonner les blocs de texte dans chaque sous-ensemble en conformité avec l'ordre défini par le chemin Hamiltonien le plus court calculé (118) ; et

 concaténer l'ordonnancement des blocs de texte dans les sous-ensembles dans un ordre final.

**FIG. 1**

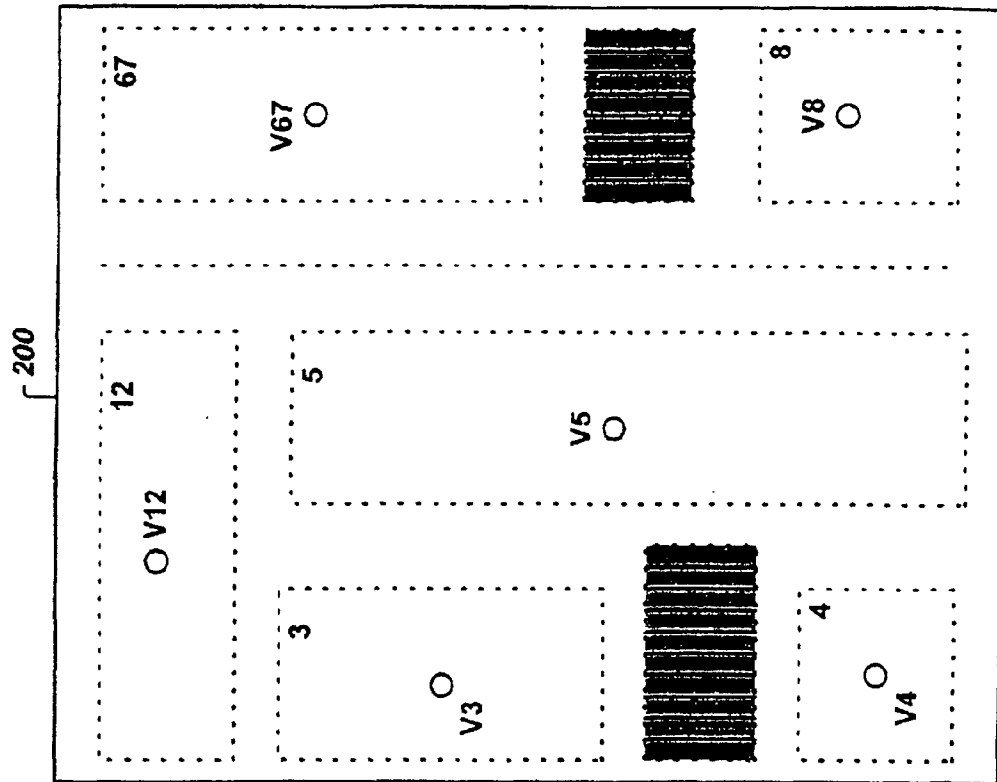


FIG. 2B

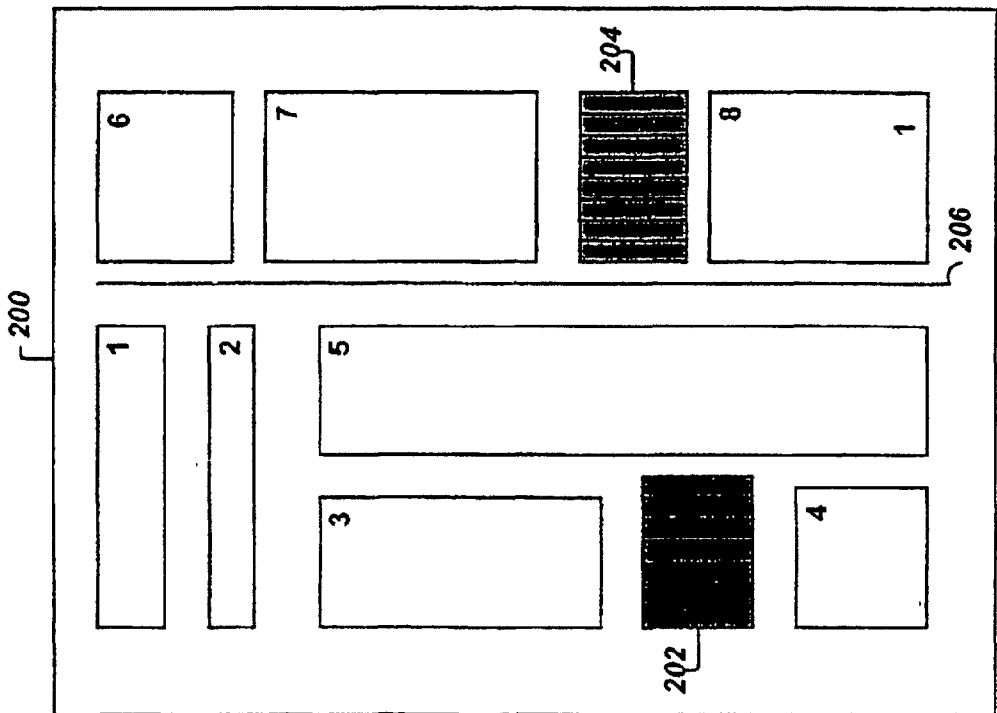
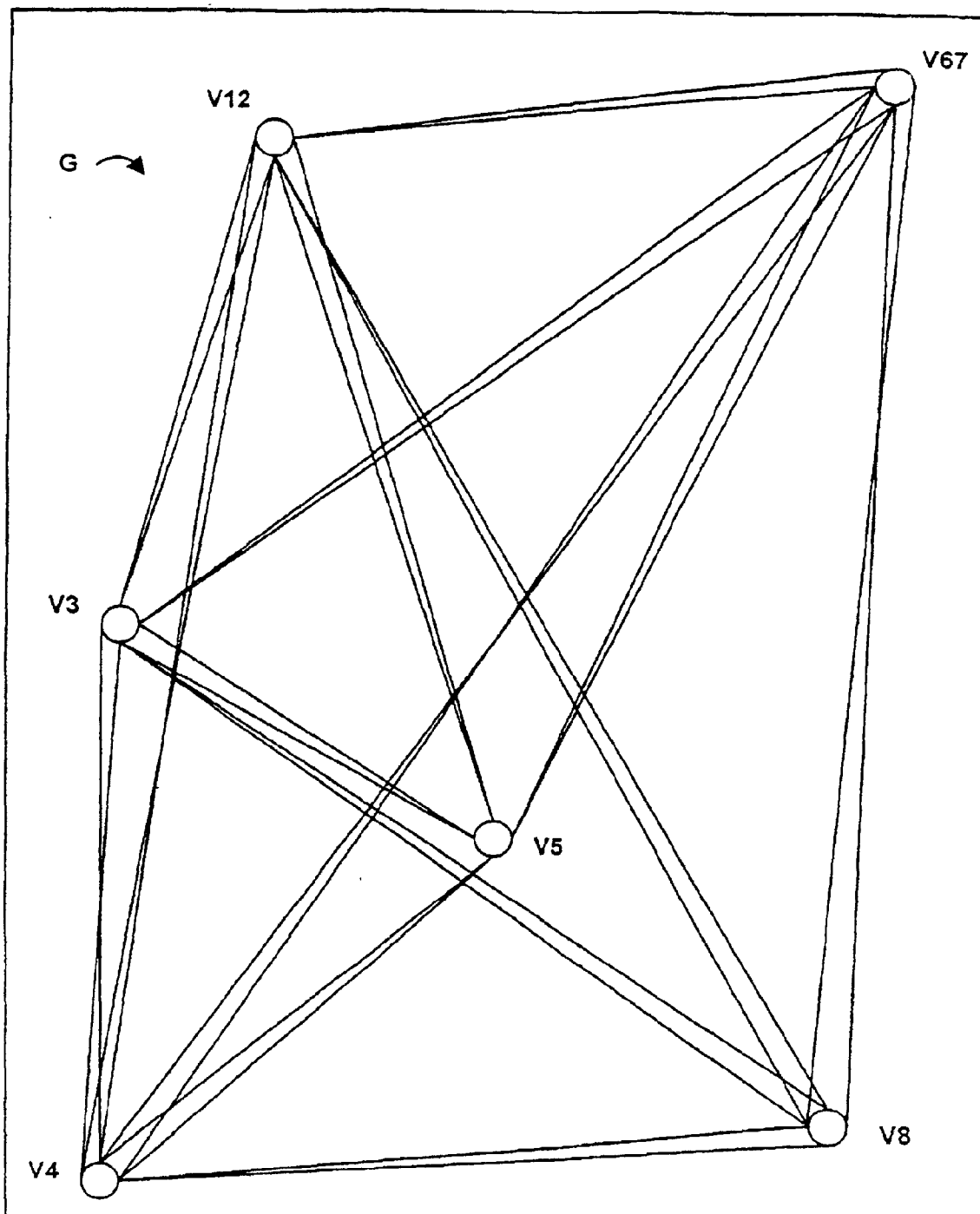


FIG. 2A

**FIG. 3**

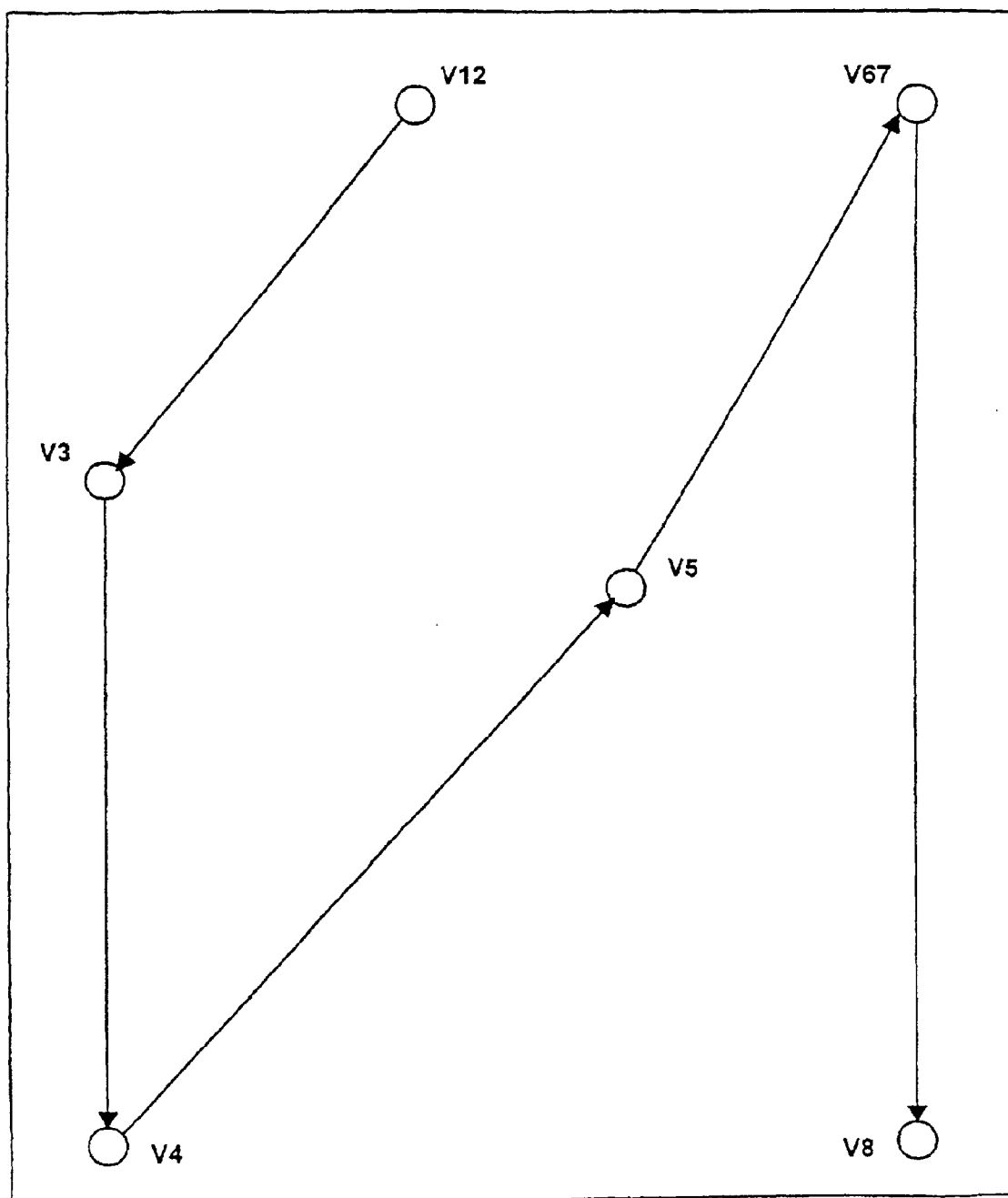
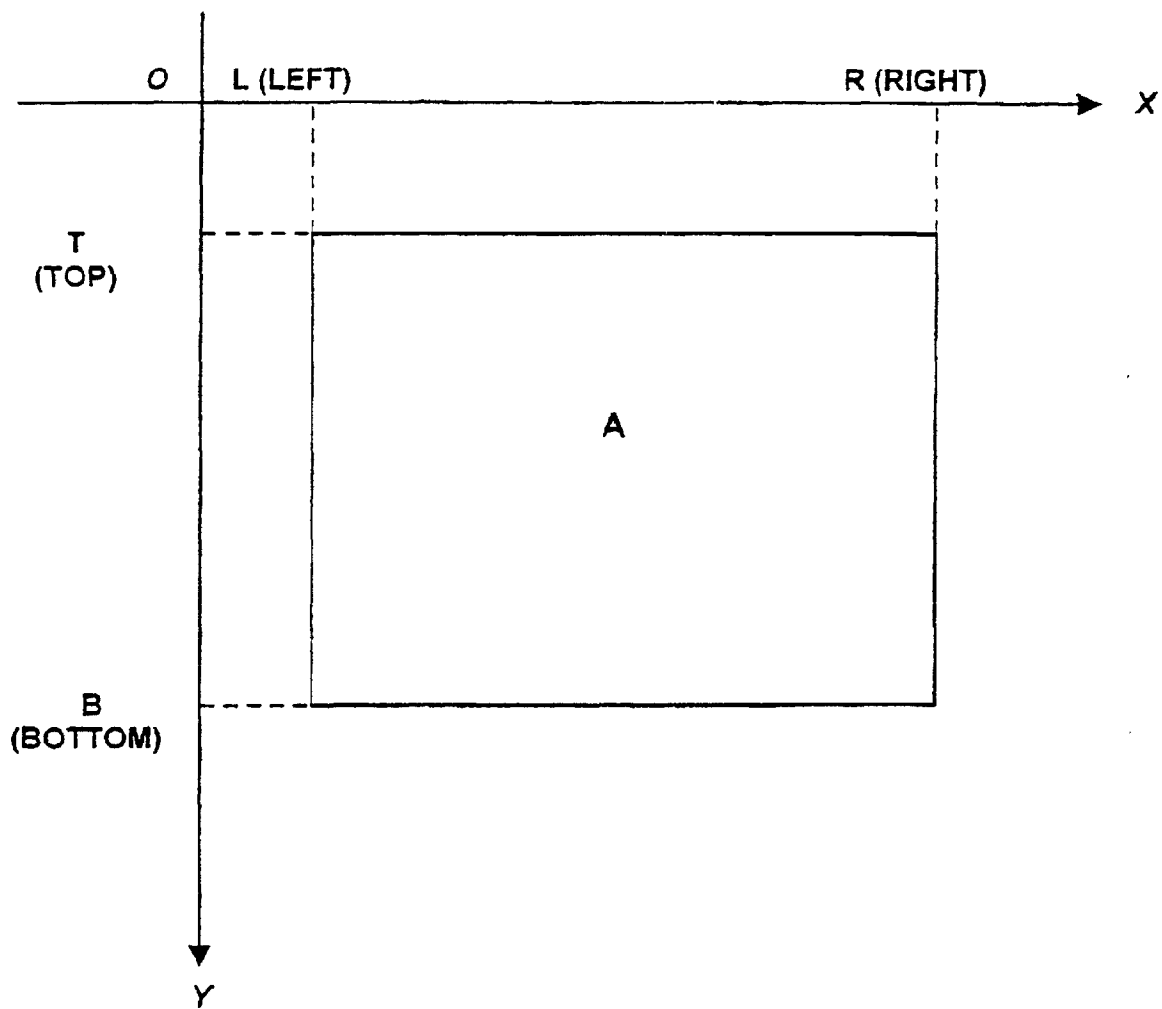
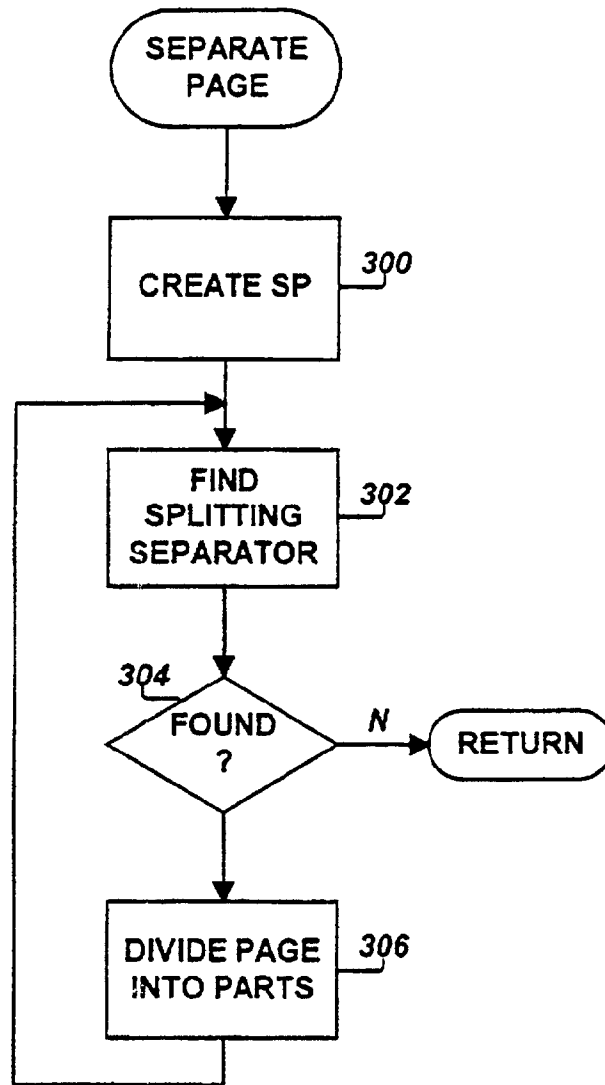


FIG. 4

**FIG. 5**

**FIG. 6**

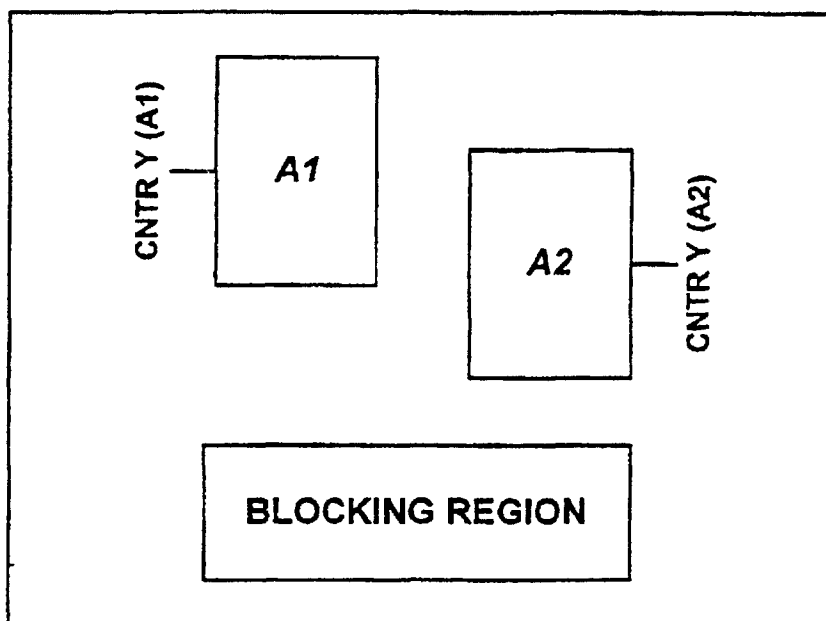


FIG. 7

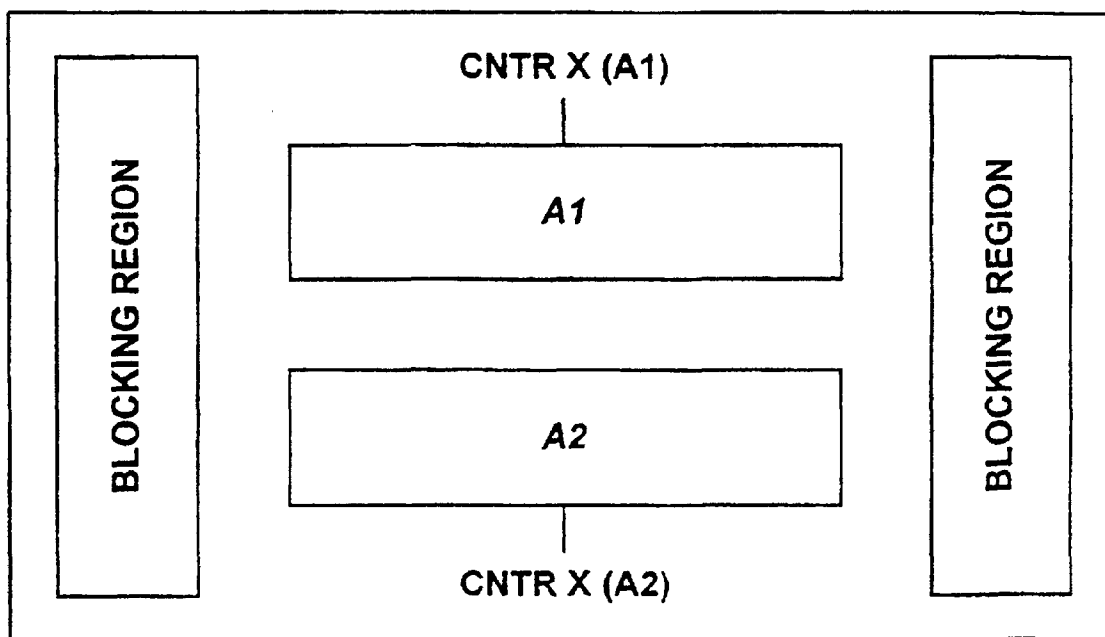


FIG. 8

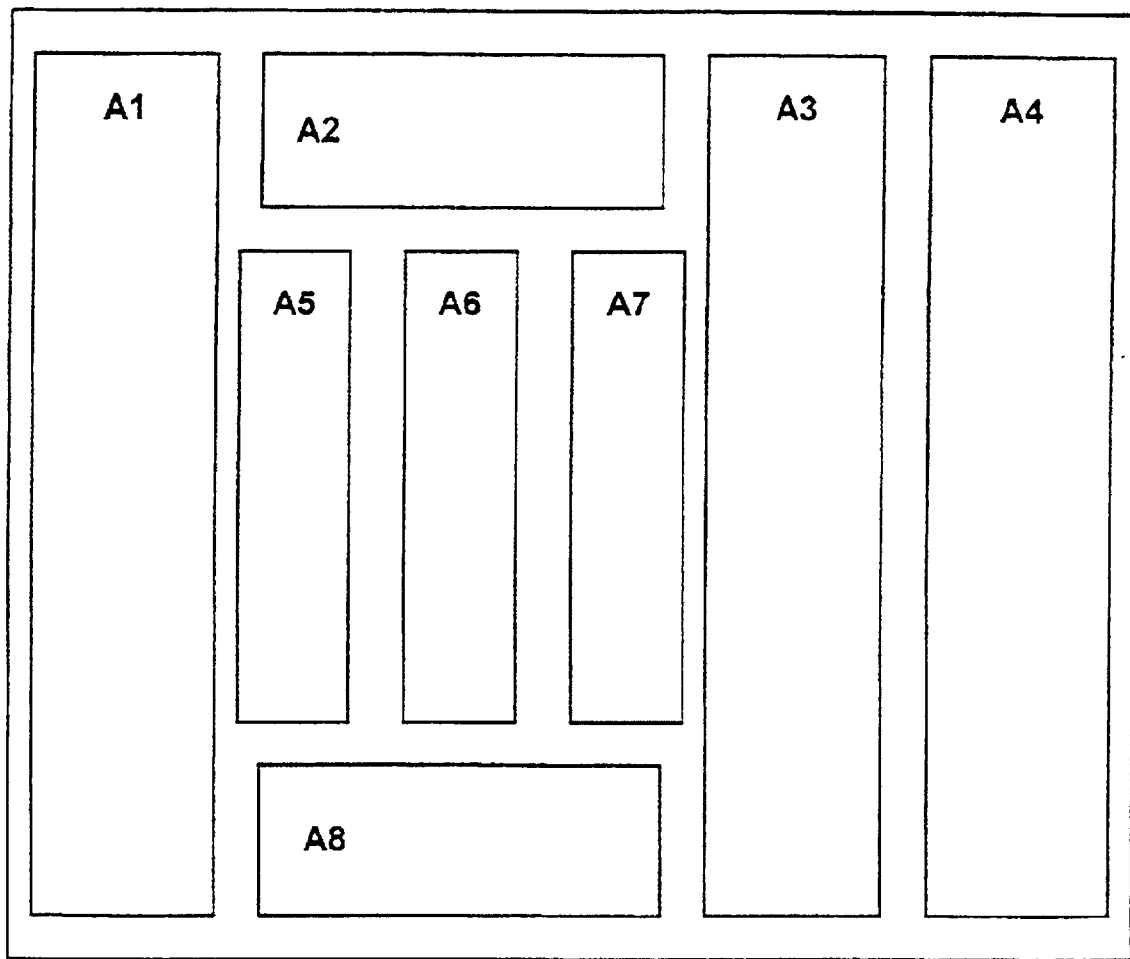
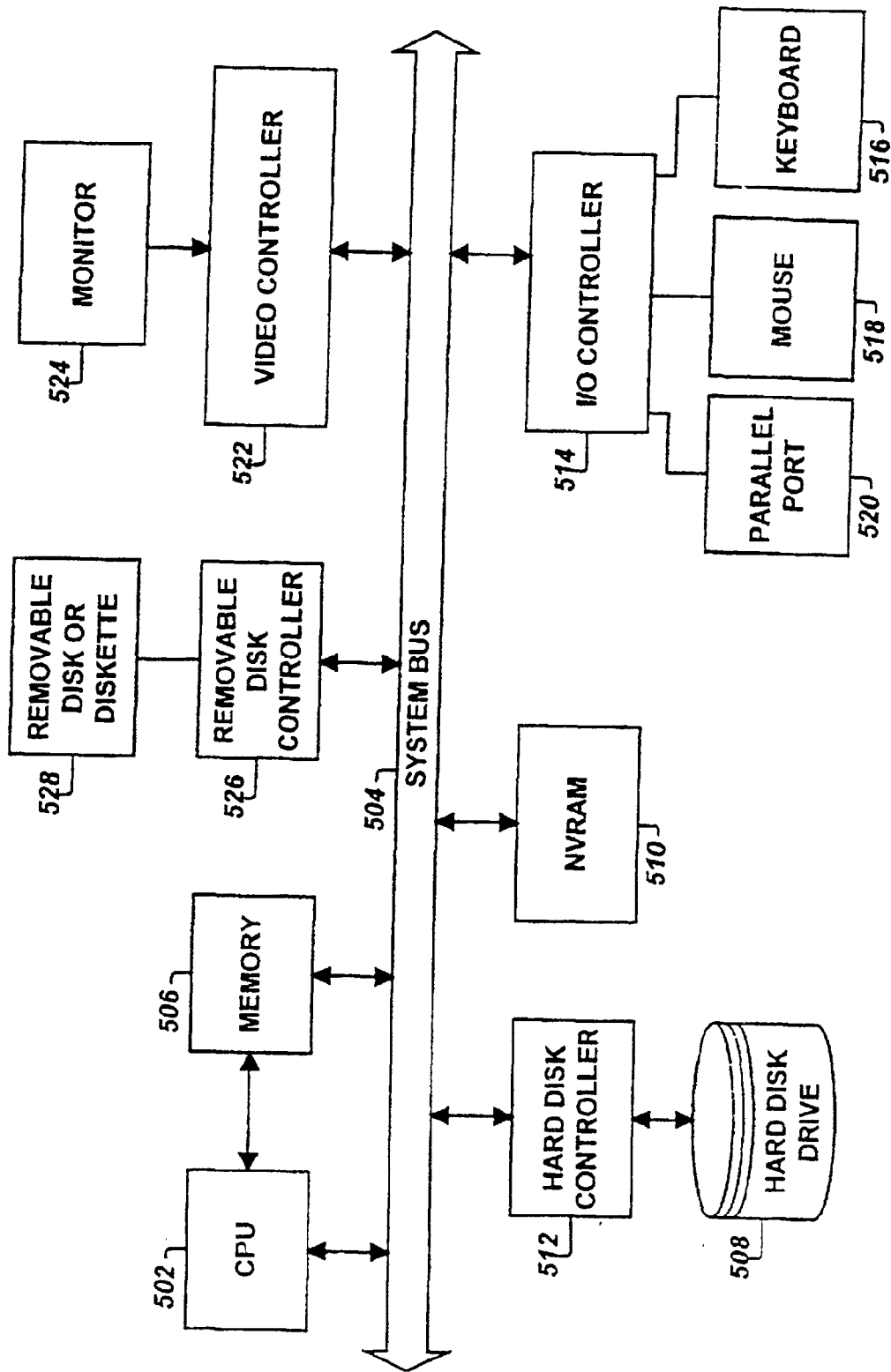


FIG. 9

**FIG. 10**